

# Support Vector Machine Classification with Indefinite Kernels

Ronny Luss\*

Alexandre d’Aspremont<sup>†</sup>

August 4, 2009

## Abstract

We propose a method for support vector machine classification using indefinite kernels. Instead of directly minimizing or stabilizing a nonconvex loss function, our algorithm simultaneously computes support vectors and a proxy kernel matrix used in forming the loss. This can be interpreted as a penalized kernel learning problem where indefinite kernel matrices are treated as noisy observations of a true Mercer kernel. Our formulation keeps the problem convex and relatively large problems can be solved efficiently using the projected gradient or analytic center cutting plane methods. We compare the performance of our technique with other methods on several standard data sets.

## 1 Introduction

Support vector machines (SVM) have become a central tool for solving binary classification problems. A critical step in support vector machine classification is choosing a suitable kernel matrix, which measures similarity between data points and must be positive semidefinite because it is formed as the Gram matrix of data points in a reproducing kernel Hilbert space. This positive semidefinite condition on kernel matrices is also known as Mercer’s condition in the machine learning literature. The classification problem then becomes a linearly constrained quadratic program. Here, we present an algorithm for SVM classification using indefinite kernels<sup>1</sup>, i.e. kernel matrices formed using similarity measures which are not positive semidefinite.

Our interest in indefinite kernels is motivated by several observations. First, certain similarity measures take advantage of application-specific structure in the data and often display excellent empirical classification performance. Unlike popular kernels used in support vector machine classification, these similarity matrices are often indefinite, so do not necessarily correspond to a reproducing kernel Hilbert space. (See Ong et al. (2004) for a discussion.)

In particular, an application of classification with indefinite kernels to image classification using Earth Mover’s Distance was discussed in Zamolotskikh and Cunningham (2004). Similarity measures for protein sequences such as the Smith-Waterman and BLAST scores are indefinite yet have provided hints for constructing useful positive semidefinite kernels such as those described in Saigo et al. (2004) or have been transformed into positive semidefinite kernels with good empirical performance (see Lanckriet et al. (2003), for example). Tangent distance similarity measures, as

---

\*ORFE Department, Princeton University, Princeton, NJ 08544. [rluss@alumni.princeton.edu](mailto:rluss@alumni.princeton.edu)

<sup>†</sup>ORFE Department, Princeton University, Princeton, NJ 08544. [aspremon@princeton.edu](mailto:aspremon@princeton.edu)

<sup>1</sup>A preliminary version of this paper appeared in the proceedings of the Neural Information Processing Systems (NIPS) 2007 conference and is available at <http://books.nips.cc/nips20.html>

described in Simard et al. (1998) or Haasdonk and Keysers (2002), are invariant to various simple image transformations and have also shown excellent performance in optical character recognition. Finally, it is sometimes impossible to prove that some kernels satisfy Mercer’s condition or the numerical complexity of evaluating the exact positive kernel is too high and a proxy (and not necessarily positive semidefinite) kernel has to be used instead (see Cuturi (2007), for example). In both cases, our method allows us to bypass these limitations. Our objective here is to derive efficient algorithms to directly use these indefinite similarity measures for classification.

Our work closely follows, in spirit, recent results on kernel learning (see Lanckriet et al. (2004) or Ong et al. (2005)), where the kernel matrix is learned as a linear combination of given kernels, and the result is explicitly constrained to be positive semidefinite. While this problem is numerically challenging, Bach et al. (2004) adapted the SMO algorithm to solve the case where the kernel is written as a positively weighted combination of other kernels. In our setting here, we never *numerically* optimize the kernel matrix because this part of the problem can be solved explicitly, which means that the complexity of our method is substantially lower than that of classical kernel learning algorithms and closer in practice to the algorithm used in Sonnenberg et al. (2006), who formulate the multiple kernel learning problem of Bach et al. (2004) as a semi-infinite linear program and solve it with a column generation technique similar to the analytic center cutting plane method we use here.

## 1.1 Current results

Several methods have been proposed for dealing with indefinite kernels in SVMs. A first direction embeds data in a pseudo-Euclidean (pE) space: Haasdonk (2005), for example, formulates the classification problem with an indefinite kernel as that of minimizing the distance between convex hulls formed from the two categories of data embedded in the pE space. The nonseparable case is handled in the same manner using reduced convex hulls. (See Bennet and Bredensteiner (2000) for a discussion on geometric interpretations in SVM.)

Another direction applies direct spectral transformations to indefinite kernels: flipping the negative eigenvalues or shifting the eigenvalues and reconstructing the kernel with the original eigenvectors in order to produce a positive semidefinite kernel (see Wu et al. (2005) and Zamolotskikh and Cunningham (2004), for example). Yet another option is to reformulate either the maximum margin problem or its dual in order to use the indefinite kernel in a convex optimization problem. One reformulation suggested in Lin and Lin (2003) replaces the indefinite kernel by the identity matrix and maintains separation using linear constraints. This method achieves good performance, but the convexification procedure is hard to interpret. Directly solving the nonconvex problem sometimes gives good results as well (see Woźnica et al. (2006) and Haasdonk (2005)) but offers no guarantees on performance.

## 1.2 Contributions

In this work, instead of directly transforming the indefinite kernel, we simultaneously learn the support vector weights and a proxy Mercer kernel matrix by penalizing the distance between this proxy kernel and the original, indefinite one. Our main result is that the kernel learning part of that problem can be solved explicitly, meaning that the classification problem with indefinite kernels can simply be formulated as a perturbation of the positive semidefinite case.

Our formulation can be interpreted as a penalized kernel learning problem with uncertainty on the input kernel matrix. In that sense, indefinite similarity matrices are seen as noisy observations of a true positive semidefinite kernel and we learn a kernel that increases the generalization performance. From a complexity standpoint, while the original SVM classification problem with indefinite kernel is nonconvex, the penalization we detail here results in a convex problem, and hence can be solved efficiently with guaranteed complexity bounds.

The paper is organized as follows. In Section 2 we formulate our main classification result and detail its interpretation as a penalized kernel learning problem. In Section 3 we describe three algorithms for solving this problem. Section 4 discusses several extensions of our main results. Finally, in Section 5, we test the numerical performance of these methods on various data sets.

### Notation

We write  $\mathbf{S}^n$  ( $\mathbf{S}_+^n$ ) to denote the set of symmetric (positive-semidefinite) matrices of size  $n$ . The vector  $e$  is the  $n$ -vector of ones. Given a matrix  $X$ ,  $\lambda_i(X)$  denotes the  $i^{th}$  eigenvalue of  $X$ .  $X_+$  is the positive part of the matrix  $X$ , i.e.  $X_+ = \sum_i \max(0, \lambda_i) v_i v_i^T$  where  $\lambda_i$  and  $v_i$  are the  $i^{th}$  eigenvalue and eigenvector of the matrix  $X$ . Given a vector  $x$ ,  $\|x\|_1 = \sum |x_i|$ .

## 2 SVM with indefinite kernels

In this section, we modify the SVM kernel learning problem and formulate a penalized kernel learning problem on indefinite kernels. We also detail how our framework applies to kernels that satisfy Mercer's condition.

### 2.1 Kernel learning

Let  $K \in \mathbf{S}^n$  be a given kernel matrix and let  $y \in \mathbf{R}^n$  be the vector of labels, with  $Y = \mathbf{diag}(y)$ , the matrix with diagonal  $y$ . We formulate the kernel learning problem as in Lanckriet et al. (2004), where the authors minimize an upper bound on the misclassification probability when using SVM with a given kernel  $K$ . This upper bound is the generalized performance measure

$$\omega_C(K) = \max_{\{0 \leq \alpha \leq C, \alpha^T y = 0\}} \alpha^T e - \mathbf{Tr}(K(Y\alpha)(Y\alpha)^T)/2 \quad (1)$$

where  $\alpha \in \mathbf{R}^n$  and  $C$  is the SVM misclassification penalty. This is also the classic 1-norm soft margin SVM problem. They show that  $\omega_C(K)$  is convex in  $K$  and solve problems of the form

$$\min_{K \in \mathcal{K}} \omega_C(K) \quad (2)$$

in order to learn an optimal kernel  $K^*$  that achieves good generalization performance. When  $\mathcal{K}$  is restricted to convex subsets of  $\mathbf{S}_+^n$  with constant trace, they show that problem (2) can be reformulated as a convex program. Further restrictions to  $\mathcal{K}$  reduce (2) to more tractable optimization problems such as semidefinite and quadratically constrained quadratic programs. Our goal is to solve a problem similar to (2) by restricting the distance between a proxy kernel used in classification and the original indefinite similarity measure.

## 2.2 Learning from indefinite kernels

The performance measure in (1) is the dual of the SVM classification problem with hinge loss and quadratic penalty. When  $K$  is positive semidefinite, this problem is a convex quadratic program. Suppose now that we are given an indefinite kernel matrix  $K_0 \in \mathbf{S}^n$ . We formulate a new instance of problem (2) by restricting  $K$  to be a positive semidefinite kernel matrix in some given neighborhood of the original (indefinite) kernel matrix  $K_0$  and solve

$$\min_{\{K \succeq 0, \|K - K_0\|_F^2 \leq \beta\}} \max_{\{\alpha^T y = 0, 0 \leq \alpha \leq C\}} \alpha^T e - \frac{1}{2} \text{Tr}(K(Y\alpha)(Y\alpha)^T)$$

in the variables  $K \in \mathbf{S}^n$  and  $\alpha \in \mathbf{R}^n$ , where the parameter  $\beta > 0$  controls the distance between the original matrix  $K_0$  and the proxy kernel  $K$ . This is the kernel learning problem (2) with  $\mathcal{K} = \{K \succeq 0, \|K - K_0\|_F^2 \leq \beta\}$ . The above problem is infeasible for small values of  $\beta$ , so we replace here the hard constraint on  $K$  by a penalty  $\rho$  on the distance between the proxy kernel and the original indefinite similarity matrix and solve instead

$$\min_{\{K \succeq 0\}} \max_{\{\alpha^T y = 0, 0 \leq \alpha \leq C\}} \alpha^T e - \frac{1}{2} \text{Tr}(K(Y\alpha)(Y\alpha)^T) + \rho \|K - K_0\|_F^2 \quad (3)$$

Because (3) is convex-concave and the inner maximization has a compact feasible set, we can switch the max and min to form the dual

$$\max_{\{\alpha^T y = 0, 0 \leq \alpha \leq C\}} \min_{\{K \succeq 0\}} \alpha^T e - \frac{1}{2} \text{Tr}(K(Y\alpha)(Y\alpha)^T) + \rho \|K - K_0\|_F^2 \quad (4)$$

in the variables  $K \in \mathbf{S}^n$  and  $\alpha \in \mathbf{R}^n$ .

We first note that problem (4) is a convex optimization problem. The inner minimization problem is a convex conic program on  $K$ . Also, as the pointwise minimum of a family of concave quadratic functions of  $\alpha$ , the solution to the inner problem is a concave function of  $\alpha$ , hence the outer optimization problem is also convex (see Boyd and Vandenberghe (2004) for further details). Thus, (4) is a concave maximization problem subject to linear constraints and is therefore a convex problem in  $\alpha$ . Our key result here is that the inner kernel learning optimization problem in (4) can be solved in closed form.

**Theorem 1** *Given a similarity matrix  $K_0 \in \mathbf{S}^n$ , a vector  $\alpha \in \mathbf{R}^n$  of support vector coefficients and the label matrix  $Y = \text{diag}(y)$ , the optimal kernel in problem (4) can be computed explicitly as:*

$$K^* = (K_0 + (Y\alpha)(Y\alpha)^T / (4\rho))_+ \quad (5)$$

where  $\rho \geq 0$  controls the penalty.

**Proof.** For a fixed  $\alpha$ , the inner minimization problem can be written out as

$$\min_{\{K \succeq 0\}} \alpha^T e + \rho (\text{Tr}(K^T K) - 2 \text{Tr}(K^T (K_0 + \frac{1}{4\rho} (Y\alpha)(Y\alpha)^T)) + \text{Tr}(K_0^T K_0))$$

where we have replaced  $\|K - K_0\|_F^2 = \text{Tr}((K - K_0)^T (K - K_0))$  and collected similar terms. Adding and subtracting the constant  $\rho \text{Tr}((K_0 + \frac{1}{4\rho} (Y\alpha)(Y\alpha)^T)^T (K_0 + \frac{1}{4\rho} (Y\alpha)(Y\alpha)^T))$  shows that the inner minimization problem is equivalent to the problem

$$\begin{aligned} & \text{minimize} && \|K - (K_0 + \frac{1}{4\rho} (Y\alpha)(Y\alpha)^T)\|_F^2 \\ & \text{subject to} && K \succeq 0 \end{aligned}$$

in the variable  $K \in \mathbf{S}^n$ , where we have dropped the remaining constants from the objective. This is the projection of the matrix  $K_0 + (Y\alpha)(Y\alpha)^T/(4\rho)$  on the cone of positive semidefinite matrices, which yields the desired result. ■

Plugging the explicit solution for the proxy kernel derived in (5) into the classification problem (4), we get

$$\max_{\{\alpha^T y = 0, 0 \leq \alpha \leq C\}} \alpha^T e - \frac{1}{2} \mathbf{Tr}(K^*(Y\alpha)(Y\alpha)^T) + \rho \|K^* - K_0\|_F^2 \quad (6)$$

in the variable  $\alpha \in \mathbf{R}^n$ , where  $(Y\alpha)(Y\alpha)^T$  is the rank one matrix with coefficients  $y_i \alpha_i \alpha_j y_j$ . Problem (6) can be cast as an eigenvalue optimization problem in the variable  $\alpha$ . Letting the eigenvalue decomposition of  $K_0 + (Y\alpha)(Y\alpha)^T/(4\rho)$  be  $VDV^T$ , we get  $K^* = VD_+V^T$ , and with  $v_i$  the  $i^{th}$  column of  $V$ , we can write

$$\begin{aligned} \mathbf{Tr}(K^*(Y\alpha)(Y\alpha)^T) &= (Y\alpha)^T V D_+ V^T (Y\alpha) \\ &= \sum_i \max \left( 0, \lambda_i \left( K_0 + \frac{1}{4\rho} (Y\alpha)(Y\alpha)^T \right) \right) (\alpha^T Y v_i)^2. \end{aligned}$$

Using the same technique, we can also rewrite the term  $\|K^* - K_0\|_F^2$  using this eigenvalue decomposition. Our original optimization problem (4) finally becomes

$$\begin{aligned} \text{maximize} \quad & \alpha^T e - \frac{1}{2} \sum_i \max(0, \lambda_i(K_0 + (Y\alpha)(Y\alpha)^T/4\rho)) (\alpha^T Y v_i)^2 \\ & + \rho \sum_i (\max(0, \lambda_i(K_0 + (Y\alpha)(Y\alpha)^T/4\rho)))^2 \\ & - 2\rho \sum_i \mathbf{Tr}((v_i v_i^T) K_0) \max(0, \lambda_i(K_0 + (Y\alpha)(Y\alpha)^T/4\rho)) + \rho \mathbf{Tr}(K_0 K_0) \\ \text{subject to} \quad & \alpha^T y = 0, 0 \leq \alpha \leq C \end{aligned} \quad (7)$$

in the variable  $\alpha \in \mathbf{R}^n$ . By construction, the objective function is concave, hence (7) is a convex optimization problem in  $\alpha$ .

A reformulation of problem (4) appears in Chen and Ye (2008) where the authors move the inner minimization problem to the constraints and get the following semi-infinite quadratically constrained linear program (SIQCLP):

$$\begin{aligned} \text{maximize} \quad & t \\ \text{subject to} \quad & \alpha^T y = 0, 0 \leq \alpha \leq C \\ & t \leq \alpha^T e - \frac{1}{2} \mathbf{Tr}(K(Y\alpha)(Y\alpha)^T) + \rho \|K - K_0\|_F^2 \quad \forall K \succeq 0. \end{aligned} \quad (8)$$

In Section 3, we describe algorithms to solve our eigenvalue optimization problem in (7), as well as an algorithm from Chen and Ye (2008) that solves the different formulation in (8), for completeness.

### 2.3 Interpretation

Our explicit solution of the optimal kernel given in (5) is the projection of a penalized rank-one update to the indefinite kernel on the cone of positive semidefinite matrices. As  $\rho$  tends to infinity, the rank-one update has less effect and in the limit, the optimal kernel is given by zeroing out the

negative eigenvalues of the indefinite kernel. This means that if the indefinite kernel contains a very small amount of noise, the best positive semidefinite kernel to use with SVM in our framework is the positive part of the indefinite kernel.

This limit as  $\rho$  tends to infinity also motivates a heuristic for transforming the kernel on the testing set. Since negative eigenvalues in the training kernel are thresholded to zero in the limit, the same transformation should occur for the test kernel. Hence, to measure generalization performance, we update the entries of the full kernel corresponding to training instances by the rank-one update resulting from the optimal solution to (7) and threshold the negative eigenvalues of the full kernel matrix to zero to produce a Mercer kernel on the test set.

## 2.4 Dual problem

As discussed above, problems (3) and (4) are dual. The inner maximization in problem (3) is a quadratic program in  $\alpha$ , whose dual is the quadratic minimization problem

$$\begin{aligned} \text{minimize} \quad & \frac{1}{2}(e - \delta + \mu + y\nu)^T(YKY)^{-1}(e - \delta + \mu + y\nu) + C\mu^T e \\ \text{subject to} \quad & \delta, \mu \geq 0. \end{aligned} \tag{9}$$

Substituting (9) for the inner maximization in problem (3) allows us to write a joint minimization problem

$$\begin{aligned} \text{minimize} \quad & \text{Tr}(K^{-1}(Y^{-1}(e - \delta + \mu + y\nu))(Y^{-1}(e - \delta + \mu + y\nu))^T)/2 + C\mu^T e + \rho\|K - K_0\|_F^2 \\ \text{subject to} \quad & K \succeq 0, \delta, \mu \geq 0 \end{aligned} \tag{10}$$

in the variables  $K \in \mathbf{S}^n$ ,  $\delta, \mu \in \mathbf{R}^n$  and  $\nu \in \mathbf{R}$ . This is a quadratic program in the variables  $\delta, \mu$  (which correspond to the constraints  $0 \leq \alpha \leq C$ ) and  $\nu$  (which is the dual variable for the constraint  $\alpha^T y = 0$ ). As we have seen earlier, any feasible solution  $\alpha \in \mathbf{R}^n$  produces a corresponding proxy kernel in (5). Plugging this kernel into problem (10) allows us to compute an upper bound on the optimum value of problem (4) by solving a simple quadratic program in the variables  $\delta, \mu, \nu$ . This result can then be used to bound the duality gap in (7) and track convergence.

## 3 Algorithms

We now detail two algorithms that can be used to solve problem (7), which maximizes a nondifferentiable concave function subject to convex constraints. An optimal point always exists since the feasible set is bounded and nonempty. For numerical stability, in both algorithms, we quadratically smooth our objective to compute a gradient. We first describe a simple projected gradient method which has numerically cheap iterations but less predictable performance in practice. We then show how to apply the analytic center cutting plane method, whose iterations are numerically more complex but which converges linearly. For completeness, we also describe an exchange method from Chen and Ye (2008) used to solve problem (8), where the numerical bottleneck is a quadratically constrained linear program solved at each iteration.

**Smoothing** Our objective contains terms of the form  $\max\{0, f(x)\}$  for some function  $f(x)$ , which are not differentiable (described in the section below). These functions are easily smoothed out by

a Moreau-Yosida regularization technique (see Hiriart-Urruty and Lemaréchal (1993), for example). We replace the max by a continuously differentiable  $\frac{\epsilon}{2}$ -approximation as follows:

$$\varphi_\epsilon(f(x)) = \max_{0 \leq u \leq 1} (uf(x) - \frac{\epsilon}{2}u^2).$$

The gradient is then given by  $\nabla \varphi_\epsilon(f(x)) = u^*(x) \nabla f(x)$  where  $u^*(x) = \operatorname{argmax} \varphi_\epsilon(f(x))$ .

**Gradient** Calculating the gradient of the objective function in (7) requires computing the eigenvalue decomposition of a matrix of the form  $X(\alpha) = K + \rho\alpha\alpha^T$ . Given a matrix  $X(\alpha)$ , the derivative of the  $i^{th}$  eigenvalue with respect to  $\alpha$  is then given by

$$\frac{\partial \lambda_i(X(\alpha))}{\partial \alpha} = v_i^T \frac{\partial X(\alpha)}{\partial \alpha} v_i \quad (11)$$

where  $v_i$  is the  $i^{th}$  eigenvector of  $X(\alpha)$ . We can then combine this expression with the smooth approximation above to obtain the gradient.

### 3.1 Computing proxy kernels

Because the proxy kernel in (5) only requires a rank one update of a (fixed) eigenvalue decomposition

$$K^* = (K_0 + (Y\alpha)(Y\alpha)^T / (4\rho))_+,$$

we now briefly recall how  $v_i$  and  $\lambda_i(X(\alpha))$  can be computed efficiently in this case (see Demmel (1997) for further details). We refer the reader to Kulis et al. (2006) for another kernel learning example using this method. Given the eigenvalue decomposition  $X = VDV^T$ , by changing basis this problem can be reduced to the decomposition of the diagonal plus rank-one matrix,  $D + \rho uu^T$ , where  $u = V^T \alpha$ . First, the updated eigenvalues are determined by solving the secular equations

$$\det(D + \rho uu^T - \lambda I) = 0,$$

which can be done in  $O(n^2)$ . While there is an explicit solution for the eigenvectors corresponding to these eigenvalues, they are not stable because the eigenvalues are approximated. This instability is circumvented by computing a vector  $\hat{u}$  such that approximate eigenvalues  $\lambda$  are the exact eigenvalues of the matrix  $D + \rho \hat{u} \hat{u}^T$ , then computing its stable eigenvectors explicitly, where both steps can be done in  $O(n^2)$  time. The key is that  $D + \rho \hat{u} \hat{u}^T$  is close enough to our original matrix so that the eigenvalues and eigenvectors are stable approximations of the true values. Finally, the eigenvectors of our original matrix are computed as  $VW$ , with  $W$  as the stable eigenvectors of  $D + \rho \hat{u} \hat{u}^T$ . Updating the eigenvalue decomposition is reduced to an  $O(n^2)$  procedure plus one matrix multiplication, which is then the complexity of one gradient computation.

We note that eigenvalues of symmetric matrices are not differentiable when some of them have multiplicities greater than one (see Overton (1992) for a discussion), but a subgradient can be used instead of the gradient in all the algorithms detailed here. Lewis (1999) shows how to compute an approximate subdifferential of the  $k$ -th largest eigenvalue of a symmetric matrix. This can then be used to form a regular subgradient of the objective function in (7) which is concave by construction.



---

**Algorithm 1** Projected gradient method

---

- 1: Compute  $\alpha_{i+1} = \alpha_i + t\nabla f(\alpha_i)$ .
  - 2: Set  $\alpha_{i+1} = p_A(\alpha_{i+1})$ .
  - 3: If  $\text{gap} \leq \epsilon$  stop, otherwise go back to step 1.
- 

### 3.2 Projected gradient method

The projected gradient method takes a steepest descent step, then projects the new point back onto the feasible region (see Bertsekas (1999), for example). We choose an initial point  $\alpha_0 \in \mathbf{R}^n$  and the algorithm proceeds as in Algorithm 1.

Here, we have assumed that the objective function is differentiable (after smoothing). The method is only efficient if the projection step is numerically cheap. The complexity of each iteration then breaks down as follows:

*Step 1.* This requires an eigenvalue decomposition that is computed in  $O(n^2)$  plus one matrix multiplication as described above. Experiments below use a stepsize of  $5/k$  for IndefiniteSVM and  $10/k$  for PerturbSVM (described in Section 4.3) where  $k$  is the iteration number. A good stepsize is crucial to performance, and must be chosen separately for each data set as there is no rule of thumb. We note that a line search would be costly here because it would require multiple eigenvalue decompositions to recalculate the objective multiple times.

*Step 2.* This is a projection onto the region  $A = \{\alpha^T y = 0, 0 \leq \alpha \leq C\}$  and can be solved explicitly by sorting the vector of entries, with cost  $O(n \log n)$ .

*Stopping Criterion.* We can compute a duality gap using the results of §2.4 where

$$K_i = (K_0 + (Y\alpha_i)(Y\alpha_i)^T / (4\rho))_+$$

is the candidate kernel at iteration  $i$  and we solve problem (1), which simply means solving a SVM problem with the positive semidefinite kernel  $K_i$ , and produces an upper bound on (7), hence a bound on the suboptimality of the current solution.

*Complexity.* The number of iterations required by this method to reach a target precision of  $\epsilon$  grows as  $O(1/\epsilon^2)$ . See Nesterov (2003) for a complete discussion.

### 3.3 Analytic center cutting plane method

The analytic center cutting plane method (ACCPM) reduces the feasible region at each iteration using a new *cut* computed by evaluating a subgradient of the objective function at the analytic center of the current feasible set, until the volume of the reduced region converges to the target precision. This method does not require differentiability. We set  $\mathcal{L}_0 = \{x \in \mathbf{R}^n \mid x^T y = 0, 0 \leq x \leq C\}$ , which we can write as  $\{x \in \mathbf{R}^n \mid A_0 x \leq b_0\}$ , to be our first localization set for the optimal solution. The method is described in Algorithm 2 (see Bertsekas (1999) for a more complete treatment of cutting plane methods).

The complexity of each iteration breaks down as follows:

*Step 1.* This step computes the analytic center of a polyhedron and can be solved in  $O(n^3)$  operations using interior point methods, for example.

*Step 2.* This simply updates the polyhedral description. It includes the gradient computation which again is  $O(n^2)$  plus one matrix multiplication.



---

**Algorithm 2** Analytic center cutting plane method

---

1: Compute  $\alpha_i$  as the analytic center of  $\mathcal{L}_i$  by solving

$$x_{i+1} = \operatorname{argmin}_{x \in \mathbf{R}^n} - \sum_{i=1}^m \log(b_i - a_i^T x)$$

where  $a_i^T$  represents the  $i^{\text{th}}$  row of coefficients from the left-hand side of  $\{x \in \mathbf{R}^n \mid A_i x \leq b_0\}$ .

2: Compute  $\nabla f(x)$  at the center  $x_{i+1}$  and update the (polyhedral) localization set

$$\mathcal{L}_{i+1} = \mathcal{L}_i \cap \{\nabla f(x_{i+1})(x - x_{i+1}) \geq 0\}$$

where  $f$  is objective in problem (7).

3: If  $m \geq 3n$ , reduce the number of constraints to  $3n$ .

4: If  $\text{gap} \leq \epsilon$  stop, otherwise go back to step 1.

---

*Step 3.* This step requires ordering the constraints according to their relevance in the localization set. One relevance measure for the  $j^{\text{th}}$  constraint at iteration  $i$  is

$$\frac{a_j^T \nabla^2 f(x_i)^{-1} a_j}{(a_j^T x_i - b_j)^2} \quad (12)$$

where  $f$  is the objective function of the analytic center problem. Computing the hessian is easy: it requires matrix multiplication of the form  $A^T D A$  where  $A$  is  $m \times n$  (matrix multiplication is kept inexpensive in this step by pruning redundant constraints) and  $D$  is diagonal. Restricting the number of constraints to  $3n$  is a rule of thumb; raising this limit increases the per iteration complexity while decreasing it increases the required number of iterations.

*Stopping Criterion.* An upper bound is computed by maximizing a first order Taylor approximation of  $f(\alpha)$  at  $\alpha_i$  over all points in an ellipsoid that covers  $\mathcal{A}_i$ , which can be computed explicitly.

*Complexity.* ACCPM is provably convergent in  $O(n(\log 1/\epsilon)^2)$  iterations when using cut elimination, which keeps the complexity of the localization set bounded. Other schemes are available with slightly different complexities: a bound of  $O(n^2/\epsilon^2)$  is achieved in Goffin and Vial (2002) using (cheaper) approximate centers, for example.

### 3.4 Exchange method for SIQCLP

The algorithm considered in Chen and Ye (2008) in order to solve problem (8) falls under a class of algorithms called exchange methods (as defined in Hettich and Kortanek (1993)). These methods iteratively solve problems constrained by a finite subset of the infinitely many constraints, where the solution at each iterate gives an improved lower bound to the maximization problem. The subproblem solved at each iteration here is

$$\begin{aligned} & \text{maximize} && t \\ & \text{subject to} && \alpha^T y = 0, 0 \leq \alpha \leq C \\ & && t \leq \alpha^T e - \frac{1}{2} \mathbf{Tr}(K_i(Y\alpha)(Y\alpha)^T) + \rho \|K_i - K_0\|_F^2 \quad i = 1, \dots, p \end{aligned} \quad (13)$$

where  $p$  is the number of constraints used to approximate the infinitely many constraints of problem (8). Let  $(t_1, \alpha_1)$  be an initial solution found by solving problem (13) with  $p = 1$  and  $K_1 = (K_0)_+$ , where  $K_0$  is the input indefinite kernel. The algorithm proceeds as in Algorithm 3 below.

---

**Algorithm 3** Exchange method

---

1: Compute  $K_{i+1}$  by solving the inner minimization problem of (4) as a function of  $\alpha_i$ .

2: Stop if

$$\alpha_i^T e - \frac{1}{2} \text{Tr}(K_{i+1}(Y\alpha_i)(Y\alpha_i)^T) + \rho \|K_{i+1} - K_0\|_F^2 \geq t_i.$$

3: Solve problem (13) with an additional constraint using  $K_{i+1}$  to get  $(t_{i+1}, \alpha_{i+1})$  and go back to step 1.

---

The complexity of each iteration breaks down as follows:

*Step 1.* This step can be solved analytically using Theorem 1. An efficient calculation of  $K_{i+1}$  can be made as in the other algorithms above using an  $O(n^2)$  procedure plus one matrix multiplication.

*Step 2 (Stopping Criterion).* The previous point  $(t_i, \alpha_i)$  is optimal if it is feasible with respect to the new constraint, in which case it is feasible for the infinitely many constraints of the original problem (8) and hence also optimal.

*Step 3.* This step requires solving a QCLP with a number of quadratic constraints equivalent to the number of iterations. As shown in Chen and Ye (2008), the QCLP can be written as a regularized version of the multiple kernel learning (MKL) problem from Lanckriet et al. (2004), where the number of constraints here is equivalent to the number of kernels in MKL. Efficient methods to solve MKL with many kernels is an active area of research, most recently in Rakotomamonjy et al. (2008). There, the authors use a gradient method to solve a reformulation of problem (13) as a smooth maximization problem. Each objective value and gradient computation requires computing a support vector machine, hence each iteration requires several SVM computations which can be speeded up using warm-starting. Furthermore, Chen and Ye (2008) prune inactive constraints at each iteration in order to decrease the number of constraints in the QCLP.

*Complexity.* No rate of convergence is known for this algorithm, but the duality gap given in Chen and Ye (2008) is shown to monotonically decrease.

### 3.5 Matlab Implementation

The first two algorithms discussed here were implemented in Matlab for the cases of indefinite (IndefiniteSVM) and positive semidefinite (PerturbSVM) kernels and can be downloaded from the authors' webpages in a package called IndefiniteSVM. The  $\rho$  penalty parameter is one-dimensional in the implementation. This package makes use of the LIBSVM code of Chang and Lin (2001) to produce suboptimality bounds and track convergence. A Matlab implementation of the exchange method (due to the authors of Chen and Ye (2008)) that uses MOSEK (MOSEK ApS 2008) to solve problem (13) is compared against the projected gradient method in Section 5.

## 4 Extensions

In this section, we extend our results to other kernel methods, namely support vector regressions and one-class support vector machines. In addition, we apply our method to using Mercer kernels

and show how to use more general penalties in our formulation.

#### 4.1 SVR with indefinite kernels

The practicality of indefinite kernels in SVM classification similarly motivates using indefinite kernels in support vector regression (SVR). We here extend the formulations in Section 2 to SVR with linear  $\epsilon$ -insensitive loss

$$\omega_C(K) = \max_{\{-C \leq \alpha \leq C, \alpha^T e = 0\}} \alpha^T y - \epsilon |\alpha| - \mathbf{Tr}(K \alpha \alpha^T) / 2 \quad (14)$$

where  $\alpha \in \mathbf{R}^n$  and  $C$  is the SVR penalty parameter. The indefinite SVR formulation follows directly as in Section 2.2 and the optimal kernel is learned by solving

$$\max_{\{\alpha^T e = 0, -C \leq \alpha \leq C\}} \min_{\{K \succeq 0\}} \alpha^T y - \epsilon |\alpha| - \frac{1}{2} \mathbf{Tr}(K \alpha \alpha^T) + \rho \|K - K_0\|_F^2 \quad (15)$$

in the variables  $K \in \mathbf{S}^n$  and  $\alpha \in \mathbf{R}^n$ , where the parameter  $\rho > 0$  controls the magnitude of the penalty on the distance between  $K$  and  $K_0$ . The following corollary to Theorem 1 provides the solution to the inner minimization problem in (15)

**Corollary 2** *Given a similarity matrix  $K_0 \in \mathbf{S}^n$  and a vector  $\alpha \in \mathbf{R}^n$  of support vector coefficients, the optimal kernel in problem (15) can be computed explicitly as*

$$K^* = (K_0 + \alpha \alpha^T / (4\rho))_+ \quad (16)$$

where  $\rho \geq 0$  controls the penalty.

The proof follows directly as in Theorem 1; the slight difference is that the vector of labels  $y$  does not appear in the optimal kernel. Plugging in (16) into (15), the resulting formulation can be rewritten as the convex eigenvalue optimization problem

$$\begin{aligned} \text{maximize} \quad & \alpha^T y - \epsilon |\alpha| - \frac{1}{2} \sum_i \max(0, \lambda_i(K_0 + \alpha \alpha^T / (4\rho))) (\alpha^T v_i)^2 \\ & + \rho \sum_i (\max(0, \lambda_i(K_0 + \alpha \alpha^T / (4\rho))))^2 \\ & - 2\rho \sum_i \mathbf{Tr}((v_i v_i^T) K_0) \max(0, \lambda_i(K_0 + \alpha \alpha^T / (4\rho))) + \rho \mathbf{Tr}(K_0 K_0) \\ \text{subject to} \quad & \alpha^T e = 0, -C \leq \alpha \leq C \end{aligned} \quad (17)$$

in the variable  $\alpha \in \mathbf{R}^n$ . Again, a proxy kernel given by (16) can be produced from any feasible solution  $\alpha \in \mathbf{R}^n$ . Plugging the proxy kernel into problem (15) allows us to compute an upper bound on the optimum value of problem (15) by solving a support vector regression problem.

#### 4.2 One-class SVM with indefinite kernels

The same reformulation can also be applied to one-class support vector machines which have the formulation (see Schölkopf and Smola (2002))

$$\omega_\nu(K) = \max_{\{0 \leq \alpha \leq \frac{1}{\nu}, \alpha^T e = 1\}} - \mathbf{Tr}(K \alpha \alpha^T) / 2 \quad (18)$$

where  $\alpha \in \mathbf{R}^n$ ,  $\nu$  is the one-class SVM parameter, and  $l$  is the number of training points. The indefinite one-class SVM formulation follows again as done for binary SVM and SVR; the optimal kernel is learned by solving

$$\max_{\{\alpha^T e = 1, 0 \leq \alpha \leq \frac{1}{\nu l}\}} \min_{\{K \succeq 0\}} -\frac{1}{2} \text{Tr}(K \alpha \alpha^T) + \rho \|K - K_0\|_F^2 \quad (19)$$

in the variables  $K \in \mathbf{S}^n$  and  $\alpha \in \mathbf{R}^n$ . The inner minimization problem is identical to that of indefinite SVR and the optimal kernel has the same form as given in Corollary 2. Plugging (16) into (19) gives another convex eigenvalue optimization problem

$$\begin{aligned} \text{maximize} \quad & -\frac{1}{2} \sum_i \max(0, \lambda_i(K_0 + \alpha \alpha^T / 4\rho)) (\alpha^T v_i)^2 \\ & + \rho \sum_i (\max(0, \lambda_i(K_0 + \alpha \alpha^T / (4\rho))))^2 \\ & - 2\rho \sum_i \text{Tr}((v_i v_i^T) K_0) \max(0, \lambda_i(K_0 + \alpha \alpha^T / (4\rho))) + \rho \text{Tr}(K_0 K_0) \\ \text{subject to} \quad & \alpha^T e = 1, 0 \leq \alpha \leq \frac{1}{\nu l} \end{aligned} \quad (20)$$

in the variable  $\alpha \in \mathbf{R}^n$ , which is identical to (17) without the first two terms in the objective and slightly different constraints. The algorithm follows almost directly the same as above for the indefinite SVR formulation.

### 4.3 Learning from Mercer kernels

While our central motivation is to use indefinite kernels for SVM classification, one would also like to analyze what happens when a Mercer kernel is used as input in (4). In this case, we learn another kernel that decreases the upper bound on generalization performance and produces perturbed support vectors. We can again interpret the input as a noisy kernel, and as such, one that will achieve suboptimal performance. If the input kernel is the best kernel to use (i.e. is not noisy), we will observe that our framework achieves optimal performance as  $\rho$  tends to infinity (through cross validation), otherwise we simply learn a better kernel using a finite  $\rho$ .

When the similarity measure  $K_0$  is positive semidefinite, the proxy kernel  $K^*$  in Theorem 1 simplifies to a rank-one update of  $K_0$

$$K^* = K_0 + (Y \alpha^*)(Y \alpha^*)^T / (4\rho) \quad (21)$$

whereas, for indefinite  $K_0$ , the solution was to project this matrix on the cone of positive semidefinite matrices. Plugging (21) into problem (4) gives:

$$\max_{\{\alpha^T y = 0, 0 \leq \alpha \leq C\}} \alpha^T e - \frac{1}{2} \text{Tr}(K_0 (Y \alpha)(Y \alpha)^T) - \frac{1}{16\rho} \sum_{i,j} (\alpha_i \alpha_j)^2, \quad (22)$$

which is the classic SVM problem given in (1) with a fourth order penalty on the support vectors. For testing in this framework, we do not need to transform the kernel, only the support vectors are perturbed. In this case, computing the gradient no longer requires eigenvalue decompositions at each iteration. Experimental results are shown in Section 5.

#### 4.4 Componentwise penalties

Indefinite SVM can be generalized further with componentwise penalties on the distance between the proxy kernel and the indefinite kernel  $K_0$ . We generalize problem (4) to

$$\max_{\{\alpha^T y = 0, 0 \leq \alpha \leq C\}} \min_{\{K \succeq 0\}} \alpha^T e - \frac{1}{2} \text{Tr}(K(Y\alpha)(Y\alpha)^T) + \sum_{i,j} H_{ij}(K_{ij} - K_{0ij})^2 \quad (23)$$

where  $H$  is now a matrix of varying penalties on the componentwise distances. For a specific class of penalties, the optimal kernel  $K^*$  can be derived explicitly as follows.

**Theorem 3** *Given a similarity matrix  $K_0 \in S^n$ , a vector  $\alpha \in \mathbf{R}^n$  of support vector coefficients and the label matrix  $Y = \text{diag}(y)$ , when  $H$  is rank-one with  $H_{ij} = h_i h_j$ , the optimal kernel in problem (23) has the explicit form*

$$K^* = W^{-1/2}((W^{1/2}(K_0 + \frac{1}{4}(W^{-1}Y\alpha^*)(W^{-1}Y\alpha^*)^T)W^{1/2})_+)W^{-1/2} \quad (24)$$

where  $W$  is the diagonal matrix with  $W_{ii} = h_i$ .

**Proof.** The inner minimization problem to problem (23) can be written out as

$$\min_{\{K \succeq 0\}} \sum_{i,j} H_{ij}(K_{ij}^2 - 2K_{ij}K_{0ij} + K_{0ij}^2) - \frac{1}{2} \sum_{i,j} y_i y_j \alpha_i \alpha_j K_{i,j}.$$

Adding and subtracting  $\sum_{i,j} H_{ij}(K_{0ij} + \frac{1}{4H_{ij}}y_i y_j \alpha_i \alpha_j)^2$ , combining similar terms, and removing remaining constants gives

$$\begin{aligned} & \text{minimize} \quad \|H^{1/2} \circ (K - (K_0 + \frac{1}{4H} \circ (Y\alpha)(Y\alpha)^T))\|_F^2 \\ & \text{subject to} \quad K \succeq 0 \end{aligned}$$

where  $\circ$  denotes the Hardamard product,  $(A \circ B)_{ij} = a_{ij}b_{ij}$ ,  $(H^{1/2})_{ij} = H_{ij}^{1/2}$ , and  $(\frac{1}{4H})_{ij} = \frac{1}{4H_{ij}}$ . This is a weighted projection problem where  $H_{ij}$  is the penalty on  $(K_{ij} - K_{0ij})^2$ . Since  $H$  is rank-one, the result follows from Theorem 3.2 of Higham (2002). ■

Notice that Theorem 3 is a generalization of Theorem 1 where we had  $H = ee^T$ . In constructing a rank-one penalty matrix  $H$ , we simply assign penalties to each training point. The componentwise penalty formulation can also be extended to true kernels. If  $K_0 \succeq 0$ , then  $K^*$  in Theorem 3 simplifies to a rank-one update of  $K_0$ :

$$K^* = K_0 + \frac{1}{4}(W^{-1/2}Y\alpha)(W^{-1/2}Y\alpha)^T \quad (25)$$

where no projection is required.

## 5 Experiments

In this section we compare the generalization performance of our technique to other methods applying SVM classification to indefinite similarity measures. We also examine classification performance using Mercer kernels. We conclude with experiments showing convergence of our algorithms. All experiments on Mercer kernels use the LIBSVM library.

## 5.1 Generalization with indefinite kernels

We compare our method for SVM classification with indefinite kernels to several kernel preprocessing techniques discussed earlier. The first three techniques perform spectral transformations on the indefinite kernel. The first, called *denoise* here, thresholds the negative eigenvalues to zero. The second transformation, called *flip*, takes the absolute value of all eigenvalues. The last transformation, *shift*, adds a constant to each eigenvalue, making them all positive. See Wu et al. (2005) for further details. We also implemented an SVM modification (denoted *Mod SVM*) suggested in Lin and Lin (2003) where a nonconvex quadratic objective function is made convex by replacing the indefinite kernel with the identity matrix. The kernel only appears in linear inequality constraints that separate the data. Finally, we compare our results with a direct use of SVM classification on the original indefinite kernel (SVM converges but the solution is only a stationary point and not guaranteed to be optimal).

We first experiment on data from the USPS handwritten digits database Hull (1994) using the indefinite Simpson score and the one-sided tangent distance kernel to compare two digits. The tangent distance is a transformation invariant measure—it assigns high similarity between an image and slightly rotated or shifted instances—and is known to perform very well on this data set. Our experiments symmetrize the one-sided tangent distance using the square of the mean tangent distance defined in Haasdonk and Keysers (2002) and make it a similarity measure by negative exponentiation. We also consider the Simpson score for this task, which is much cheaper to compute (a ratio comparing binary pixels). We finally analyze three data sets (diabetes, german and ala) from the UCI repository (Asuncion and Newman 2007) using the indefinite sigmoid kernel.

The data is randomly divided into training and testing data. We apply 5-fold cross validation and use an average of the accuracy and recall measures (described below) to determine the optimal parameters  $C$ ,  $\rho$ , and any kernel inputs. We then train a model with the full training set and optimal parameters and test on the independent test set.

Table 1 provides summary statistics for these data sets, including the minimum and maximum eigenvalues of the training similarity matrices. We observe that the Simpson are highly indefinite, while the one-sided tangent distance kernel is nearly positive semidefinite. The spectrum of sigmoid kernels varies greatly across examples because it is very sensitive to the sigmoid kernel parameters. Table 2 compares accuracy, recall, and their average for *denoise*, *flip*, *shift*, modified SVM, direct SVM and the indefinite SVM algorithm described in this work.

Based on the interpretation from Section 2.3, Indefinite SVM should be expected to perform at least as well as *denoise*; if *denoise* were a good transformation, then cross-validation over  $\rho$  should choose a high penalty that makes Indefinite SVM and *denoise* nearly equivalent. The rank-one update provides more flexibility for the transformation and similarities concerning data points  $x_i$  that are easily classified ( $\alpha_i = 0$ ) are not modified by the rank-one update. Further interpretation for the specific rank-one update is not currently known. However, Chen et al. (2009) recently proposed spectrum modifications in a similar manner to Indefinite SVM. Rather than perturb the entire indefinite similarity matrix, they perturb the spectrum directly allowing improvements over the *denoise* as well as *flip* transformations. They also note that Indefinite SVM might perform better on sparse kernels because the rank-one update may then allow inference of hidden relationships.

We observe that Indefinite SVM performs comparably on all USPS examples (slightly better for the highly indefinite Simpson kernels), which are relatively easy classification problems. As expected, classification using the tangent distance outperforms classification with the Simpson score but, as mentioned above, the Simpson score is cheaper to compute. We also note that other

Data Set	# Train	# Test	$\lambda_{min}$	$\lambda_{max}$
USPS-3-5-SS	767	773	-70.00	903.94
USPS-3-5-TD1	767	773	-0.31	764.72
USPS-4-6-SS	829	857	-74.38	819.36
USPS-4-6-TD1	829	857	-0.72	771.07
diabetes-sig	384	384	-.65	211.62
german-sig	500	500	-928.10	8.50
a1a-sig	803	802	-.01	84.44

**Table 1:** Summary statistics for the various data sets used in our experiments. The USPS data comes from the USPS handwritten digits database, the other data sets are taken from the UCI repository. *SS* refers to the Simpson kernel, *TD1* to the one-sided tangent distance kernel, and *sig* to the sigmoid kernel. Training and testing sets were divided randomly. Notice that the Simpson kernels are mostly highly indefinite while the one-sided tangent distance kernel is nearly positive semidefinite. The sigmoid kernel is highly indefinite depending on the parametrization. Statistics for sigmoid kernels refer to the optimal kernel parameterized under cross validation with Indefinite SVM. Spectrums are based on the full kernel, i.e. combining training and testing data.

documented classification results on this USPS data set perform multi-classification, while here we only perform binary classification. Classification of the UCI data sets with sigmoid kernels is more difficult (as demonstrated by lower performance measures). Indefinite SVM here is the only technique that outperforms in at least one of the measures across all three data sets.

## 5.2 Generalization with Mercer kernels

Using this time linear and gaussian (both positive semidefinite, i.e. Mercer) kernels on the USPS data set, we now compare classification performance using regular SVM and the penalized kernel learning problem (22) of Section 4.3, which we call PerturbSVM here. We also test these two techniques on positive semidefinite kernels formed using noisy USPS data sets (created by adding uniformly distributed noise in  $[-1,1]$  to each pixel before normalizing to  $[0,1]$ ), in which case PerturbSVM can be seen as optimally denoised support vector machine classification. We again cross-validate on a training set and test on the same independent group of examples used in the experiments above. Optimal parameters from classification of unperturbed data were used to train classifiers for perturbed data. Results are summarized in Table 3.

These results show that PerturbSVM performs at least as well in almost all cases. As expected, noise decreased generalization performance in all experiments. Except in the *USPS-4-6-gaussian* example, the value of  $\rho$  selected was not the highest possible for each test where PerturbSVM outperforms SVM in at least one measure; this implies that the support vectors were perturbed to improve classification. Overall, when zero or moderate noise is present, PerturbSVM does improve performance over regular SVM as shown. When too much noise is present however (for example, pixel data with range in  $[-1,1]$  was modified with uniform noise in  $[-2,2]$  before being normalized to  $[0,1]$ ), the performance of both techniques is comparable.



Data Set	Measure	Denoise	Flip	Shift	Mod SVM	SVM	Indefinite SVM
USPS-3-5-SS	Accuracy	95.47	95.21	93.27	<b>96.12</b>	69.47	95.73
	Recall	94.50	94.50	94.98	96.17	67.94	<b>97.13</b>
	Average	94.98	94.86	94.12	96.15	68.71	<b>96.43</b>
USPS-3-5-TD1	Accuracy	<b>98.58</b>	98.45	<b>98.58</b>	98.19	<b>98.58</b>	98.45
	Recall	<b>98.56</b>	98.33	<b>98.56</b>	97.85	<b>98.56</b>	98.33
	Average	<b>98.57</b>	98.39	<b>98.57</b>	98.02	<b>98.57</b>	98.39
USPS-4-6-SS	Accuracy	<b>98.60</b>	98.25	96.73	<b>98.60</b>	84.36	98.25
	Recall	99.32	99.32	96.61	99.32	81.72	<b>99.77</b>
	Average	98.96	98.79	96.67	98.96	83.04	<b>99.01</b>
USPS-4-6-TD1	Accuracy	<b>99.30</b>	<b>99.30</b>	99.18	99.18	<b>99.30</b>	<b>99.30</b>
	Recall	<b>99.77</b>	<b>99.77</b>	99.55	99.55	<b>99.77</b>	<b>99.77</b>
	Average	<b>99.54</b>	<b>99.54</b>	99.37	99.37	<b>99.54</b>	<b>99.54</b>
diabetes-sig	Accuracy	74.48	74.74	76.56	76.04	73.70	<b>77.08</b>
	Recall	78.40	76.80	<b>89.60</b>	78.40	76.40	89.20
	Average	76.44	75.77	83.08	77.22	75.05	<b>83.14</b>
german-sig	Accuracy	70.40	70.40	<b>75.60</b>	72.60	69.40	62.80
	Recall	78.00	78.00	46.67	66.00	80.00	<b>85.33</b>
	Average	74.20	74.20	61.13	69.30	<b>74.70</b>	74.07
ala-sig	Accuracy	74.06	76.18	75.69	78.55	75.69	<b>82.92</b>
	Recall	87.31	87.82	87.31	<b>89.34</b>	87.82	81.73
	Average	80.69	82.00	81.50	<b>83.95</b>	81.75	82.32

**Table 2:** Indefinite SVM performs favorably for the highly indefinite Simpson kernels. Performance is comparable for the nearly positive semidefinite one-sided tangent distance kernel. Comparable performance with sigmoid kernels is more consistent with indefinite SVM across data sets. The performance measures are: Accuracy =  $\frac{TP+TN}{TP+TN+FP+FN}$ , Recall =  $\frac{TP}{TP+FN}$ , and Average = (Accuracy + Recall)/2.

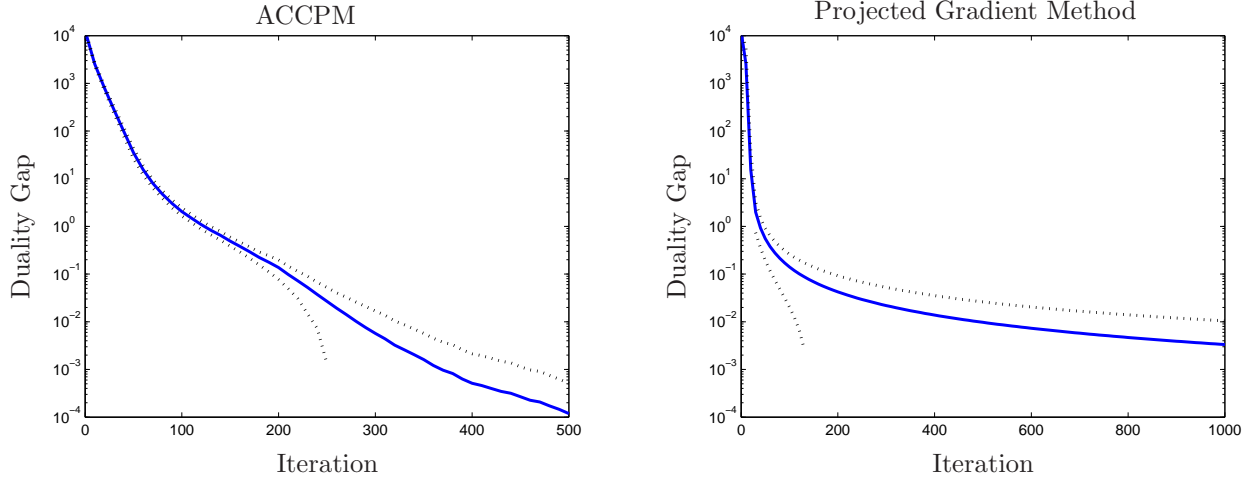
### 5.3 Convergence

We ran our two algorithms on data sets created by randomly perturbing the four USPS data sets used above. Average results and standard deviation are displayed in Figure 1 in semilog scale (note that the codes were not stopped here and that the target duality gap improvement is usually much smaller than  $10^{-8}$ ). As expected, ACCPM converges much faster (in fact linearly) to a higher precision, while each iteration requires solving a linear program of size  $n$ . The gradient projection method converges faster in the beginning but stalls at higher precision, however each iteration only requires a rank one update on an eigenvalue decomposition.

We finally examine the computing time of IndefiniteSVM using the projected gradient method

Data Set	Measure	Unperturbed		Noisy	
		SVM	Perturb SVM	SVM	Perturb SVM
USPS-3-5-linear	Accuracy	<b>96.25</b>	96.12	90.27	<b>93.16</b>
	Recall	95.69	<b>95.93</b>	90.00	<b>92.87</b>
	Average	95.97	<b>96.03</b>	90.14	<b>93.01</b>
USPS-4-6-linear	Accuracy	<b>99.07</b>	<b>99.07</b>	97.39	<b>97.97</b>
	Recall	99.10	<b>99.32</b>	97.34	<b>98.13</b>
	Average	99.08	<b>99.19</b>	97.36	<b>98.05</b>
USPS-3-5-gaussian	Accuracy	<b>97.67</b>	97.54	92.11	<b>93.57</b>
	Recall	<b>98.09</b>	97.37	91.27	<b>92.89</b>
	Average	<b>97.88</b>	97.46	91.69	<b>93.23</b>
USPS-4-6-gaussian	Accuracy	99.18	<b>99.30</b>	<b>98.00</b>	97.99
	Recall	<b>99.55</b>	<b>99.55</b>	98.15	<b>98.19</b>
	Average	99.37	<b>99.42</b>	98.08	<b>98.09</b>

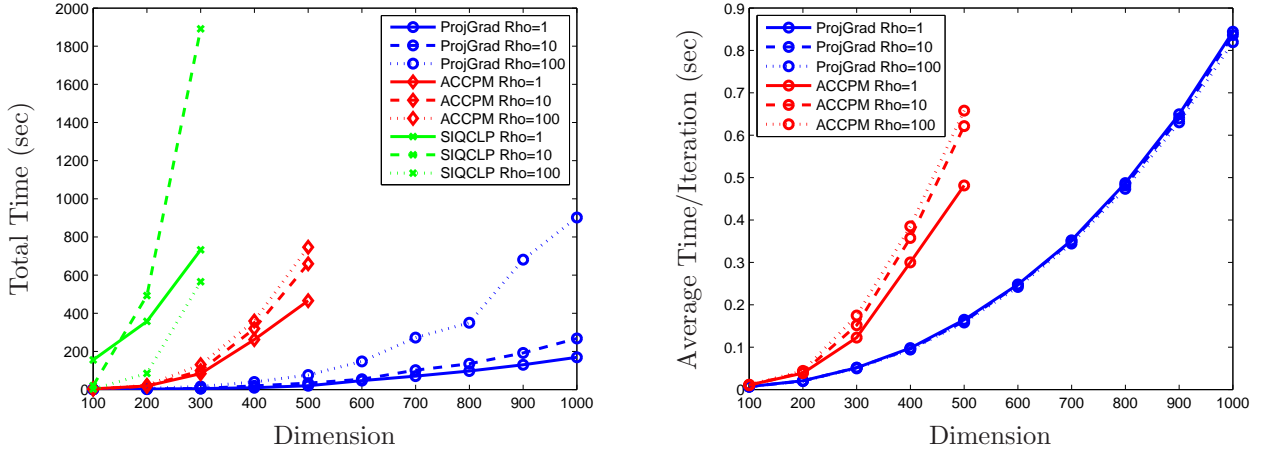
**Table 3:** Performance measures for USPS data using linear and gaussian kernels. *Unperturbed* refers to classification of the original data and *Noisy* refers to classification of data that is perturbed by uniform noise. Perturb SVM perturbs the support vectors to improve generalization. However, performance is lower for both techniques in the presence of high noise.



**Figure 1:** Convergence plots for ACCPM (left) and projected gradient method (right) on random subsets of the USPS-SS-3-5 data set (average gap versus iteration number, dashed lines at plus and minus one standard deviation). ACCPM converges linearly to a higher precision while the gradient projection method converges faster in the beginning but stalls at a higher precision.

and ACCPM and compare them with the SIQCLP method of Chen and Ye (2008). Figure 2 shows total runtime (left) and average iteration runtime (right) for varying problem dimensions on an example from the USPS data with Simpson kernel. Experiments are averaged over 10 random data subsets and we fix  $C = 10$  with a tolerance of .1 for the duality gap. For the projected gradient method, increasing  $\rho$  increases the number of iterations to converge; notice that the average time per iteration does not vary over  $\rho$ . SIQCLP also requires more iterations to converge for higher  $\rho$ , however the average iteration time seems to be less for higher  $\rho$ , so no clear pattern is seen when varying  $\rho$ . Note that the number of iterations required varies widely (between 100 and 2000 iterations in this experiment) as a function of  $\rho$ ,  $C$ , the chosen kernel and the stepsize.

Results for ACCPM and SIQCLP are shown only up to dimensions 500 and 300, respectively, because this sufficiently demonstrates that the projected gradient method is more efficient. ACCPM clearly suffers from the complexity of the analytic center problem each iteration. However, improvements can be made in the SIQCLP implementation such as using a regularized version of an efficient MKL solver (e.g. Rakotomamonjy et al. (2008)) to solve problem (13) rather than MOSEK. SIQCLP is also useful because it makes a connection between the indefinite SVM formulation and multiple kernel learning. We observed from experiments that the duality gap found from SIQCLP is tighter than the upper bound on the duality gap used for the projected gradient method. This could potentially be used to create a better stopping condition, however the complexity to derive the tighter duality gap (solving regularized MKL) is much higher than that to compute our current gap (solving a single SVM).



**Figure 2:** Total time versus dimension (left) and average time per iteration versus dimension (right) using projected gradient and ACCPM IndefiniteSVM and SIQCLP (only for total time). The number of iterations for convergence varies from 100 for the smallest dimension to 2000 for the largest dimension in this example which uses a Simpson kernel on the USPS 3-5 data.

## 6 Conclusion

We have proposed a technique for support vector machine classification with indefinite kernels, using a proxy kernel which can be computed explicitly. We also show how this framework can be

used to improve generalization performance with potentially noisy Mercer kernels, as well as extend it to other kernel methods such as support vector regression and one-class support vector machines. We give two provably convergent algorithms for solving this problem on relatively large data sets. Our initial experiments show that our method fares quite favorably compared to other techniques handling indefinite kernels in the SVM framework and, in the limit, provides a clear interpretation for some of these heuristics.

## Acknowledgements

We are very grateful to Mátyás Sustik for his rank-one update eigenvalue decomposition code and to Jianhui Chen and Jieping Ye for their SIQCLP Matlab code. We would also like to acknowledge support from NSF grant DMS-0625352, NSF CDI grant SES-0835550, a NSF CAREER award, a Peek junior faculty fellowship and a Howard B. Wentz Jr. junior faculty award.

## References

- Asuncion, A. and Newman, D. (2007), *UCI Machine Learning Repository*, University of California, Irvine, School of Information and Computer Sciences. <http://www.ics.uci.edu/~mllearn/MLRepository.html>.
- Bach, F. R., Lanckriet, G. R. G. and Jordan, M. I. (2004), ‘Multiple kernel learning, conic duality, and the SMO algorithm’, *Proceedings of the 21st International Conference on Machine Learning* . 8 pages.
- Bennet, K. P. and Bredensteiner, E. J. (2000), ‘Duality and geometry in svm classifiers’, *Proceedings of the 17th International conference on Machine Learning* pp. 57–64.
- Bertsekas, D. (1999), *Nonlinear Programming, 2nd Edition*, Athena Scientific.
- Boyd, S. and Vandenberghe, L. (2004), *Convex Optimization*, Cambridge University Press.
- Chang, C.-C. and Lin, C.-J. (2001), ‘LIBSVM: a library for support vector machines’. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- Chen, J. and Ye, J. (2008), ‘Training SVM with indefinite kernels’, *Proceedings of the 25th International Conference on Machine Learning* . 8 pages.
- Chen, Y., Gupta, M. R. and Recht, B. (2009), ‘Learning kernels from indefinite similarities’, *Proceedings of the 26th International Conference on Machine Learning* . 8 pages.
- Cuturi, M. (2007), ‘Permanents, transport polytopes and positive definite kernels on histograms’, *Proceedings of the Twentieth International Joint Conference on Artificial Intelligence* pp. 732–737.
- Demmel, J. W. (1997), *Applied Numerical Linear Algebra*, SIAM.
- Goffin, J.-L. and Vial, J.-P. (2002), ‘Convex nondifferentiable optimization: A survey focused on the analytic center cutting plane method’, *Optimization Methods and Software* **17**(5), 805–867.
- Haasdonk, B. (2005), ‘Feature space interpretation of SVMs with indefinite kernels’, *IEEE Transactions on Pattern Analysis and Machine Intelligence* **27**(4), 482–492.
- Haasdonk, B. and Keysers, D. (2002), ‘Tangent distance kernels for support vector machines’, *Proc. of the 16th Int. Conf. on Pattern Recognition* **2**, 864–868.
- Hettich, R. and Kortanek, K. O. (1993), ‘Semi-infinite programming: Theory, methods, and applications’, *SIAM Review* **35**(3), 380–429.
- Higham, N. (2002), ‘Computing the nearest correlation matrix—a problem from finance’, *IMA Journal of Numerical Analysis* **22**, 329–343.

- Hiriart-Urruty, J.-B. and Lemaréchal, C. (1993), *Convex Analysis and Minimization Algorithms*, Springer.
- Hull, J. J. (1994), ‘A database for handwritten text recognition research’, *IEEE Transactions on Pattern Analysis and Machine Intelligence* **16**(5), 550–554.
- Kulis, B., Sustik, M. and Dhillon, I. (2006), ‘Learning low-rank kernel matrices’, *Proceedings of the 23rd International Conference on Machine Learning* pp. 505–512.
- Lanckriet, G. R. G., Cristianini, N., Bartlett, P., Ghaoui, L. E. and Jordan, M. I. (2004), ‘Learning the kernel matrix with semidefinite programming’, *Journal of Machine Learning Research* **5**, 27–72.
- Lanckriet, G. R. G., Cristianini, N., Jordan, M. I. and Noble, W. S. (2003), ‘Kernel-based integration of genomic data using semidefinite programming’. In *Kernel Methods in Computational Biology*, MIT Press.
- Lewis, A. (1999), ‘Nonsmooth analysis of eigenvalues’, *Mathematical Programming* **84**, 1–24.
- Lin, H.-T. and Lin, C.-J. (2003), ‘A study on sigmoid kernel for SVM and the training of non-PSD kernels by SMO-type methods’. National Taiwan University, Taipei, Taiwan, Department of Computer Science and Information Engineering.
- MOSEK ApS (2008), ‘The MOSEK optimization tools manual. version 5.0, revision 105.’. Software available at <http://www.mosek.com>.
- Nesterov, Y. (2003), *Introductory Lectures on Convex Optimization*, Springer.
- Ong, C. S., Mary, X., Canu, S. and Smola, A. J. (2004), ‘Learning with non-positive kernels’, *Proceedings of the 21st International Conference on Machine Learning* . 8 pages.
- Ong, C. S., Smola, A. J. and Williamson, R. C. (2005), ‘Learning the kernel with hyperkernels’, *Journal of Machine Learning Research* **6**, 1043–1071.
- Overton, M. (1992), ‘Large-scale optimization of eigenvalues’, *SIAM Journal on Optimization* **2**(1), 88–120.
- Rakotomamonjy, A., Bach, F., Canu, S. and Grandvalet, Y. (2008), ‘SimpleMKL’, *Journal of Machine Learning Research* (9), 2491–2521.
- Saigo, H., Vert, J. P., Ueda, N. and Akutsu, T. (2004), ‘Protein homology detection using string alignment kernels’, *Bioinformatics* **20**(11), 1682–1689.
- Schölkopf, B. and Smola, A. (2002), *Learning with Kernels*, The MIT Press.
- Simard, P. Y., Cun, Y. A. L., Denker, J. S. and Victorri, B. (1998), ‘Transformation invariance in pattern recognition-tangent distance and tangent propagation’, *Lecture Notes in Computer Science* **1524**, 239–274.
- Sonnenberg, S., Rätsch, G., Schäfer, C. and Schölkopf, B. (2006), ‘Large scale multiple kernel learning’, *Journal of Machine Learning Research* **7**, 1531–1565.
- Woźnica, A., Kalousis, A. and Hilario, M. (2006), ‘Distances and (indefinite) kernels for set of objects’, *Proceedings of the 6th International Conference on Data Mining* pp. 1151–1156.
- Wu, G., Chang, E. Y. and Zhang, Z. (2005), ‘An analysis of transformation on non-positive semidefinite similarity matrix for kernel machines’, *Proceedings of the 22nd International Conference on Machine Learning* . 8 pages.
- Zamolotskikh, A. and Cunningham, P. (2004), ‘An assessment of alternative strategies for constructing EMD-based kernel functions for use in an SVM for image classification’, *Technical Report UCD-CSI-2007-3* . University College Dublin, School of Computer Science and Informatics.